

# Probabilistic cell/domain-type assignment of spatial transcriptomics data with SpatialAnno

Xingjie Shi<sup>1,\*†</sup>, Yi Yang<sup>2,†</sup>, Xiaohui Ma<sup>3</sup>, Yong Zhou<sup>1</sup>, Zhenxing Guo<sup>4</sup>, Chaolong Wang<sup>5</sup> and Jin Liu<sup>4,\*</sup>

<sup>1</sup>KLATASDS-MOE, Academy of Statistics and Interdisciplinary Sciences, School of Statistics, East China Normal University, Shanghai 200062, China

<sup>2</sup>The Key Laboratory of Developmental Genes and Human Disease, School of Life Science and Technology, Southeast University, Nanjing 210018, China

<sup>3</sup>College of Life Sciences, Nanjing University, Nanjing 210033, China

<sup>4</sup>School of Data Science, The Chinese University of Hong Kong-Shenzhen, Shenzhen 518172, China

<sup>5</sup>Department of Epidemiology and Biostatistics, School of Public Health, Tongji Medical College, Huazhong University of Science and Technology, Wuhan 430070, China

\*To whom correspondence should be addressed. Email: xjshi@fem.ecnu.edu.cn

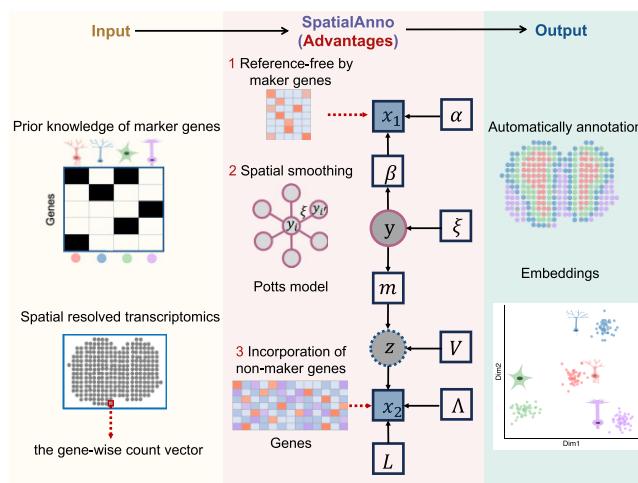
Correspondence may also be addressed to Jin Liu. Email: liujinlab@cuhk.edu.cn

†The authors wish it to be known that, in their opinion, these authors should be regarded as Joint First Authors.

## Abstract

In the analysis of both single-cell RNA sequencing (scRNA-seq) and spatially resolved transcriptomics (SRT) data, classifying cells/spots into cell/domain types is an essential analytic step for many secondary analyses. Most of the existing annotation methods have been developed for scRNA-seq datasets without any consideration of spatial information. Here, we present SpatialAnno, an efficient and accurate annotation method for spatial transcriptomics datasets, with the capability to effectively leverage a large number of non-marker genes as well as ‘qualitative’ information about marker genes without using a reference dataset. Uniquely, SpatialAnno estimates low-dimensional embeddings for a large number of non-marker genes via a factor model while promoting spatial smoothness among neighboring spots via a Potts model. Using both simulated and four real spatial transcriptomics datasets from the 10x Visium, ST, Slide-seqV1/2, and seqFISH platforms, we showcase the method’s improved spatial annotation accuracy, including its robustness to the inclusion of marker genes for irrelevant cell/domain types and to various degrees of marker gene misspecification. SpatialAnno is computationally scalable and applicable to SRT datasets from different platforms. Furthermore, the estimated embeddings for cellular biological effects facilitate many downstream analyses.

## Graphical abstract



## Introduction

With the rapid advancement of spatially resolved transcriptomics (SRT) technologies, it has become feasible to comprehensively characterize the gene expression profiles of tis-

sues while retaining information on their physical locations. Among the already developed SRT methods, *in situ* hybridization (ISH) technologies, such as MERFISH (1) and seqFISH (2), provide single-molecule resolution for targeted genes but

Received: January 25, 2023. Revised: October 3, 2023. Editorial Decision: October 7, 2023. Accepted: October 20, 2023

© The Author(s) 2023. Published by Oxford University Press on behalf of Nucleic Acids Research.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

require prior knowledge of the genes of interest. To enable single-cell analysis, cell segmentation must be performed to assign transcripts to individual cells. Alternatively, *in situ* capturing technologies, such as 10x Visium, Slide-seqV1/2 (3) and Stereo-seq (4), are unbiased and provide transcriptome-wide expression measurements. Among the *in situ* capturing technologies, there has been a dramatic improvement in spatial resolution, with spot sizes ranging from 55 µm in 10x Visium, 10 µm in Slide-seqV2, to <1 µm in Stereo-seq. These SRT technologies provide an opportunity to study how the spatial organization of gene expression in tissues relates to tissue functions (5). To characterize the transcriptomic landscape within a spatial context, assigning cell/domain types in relation to tissue location is an essential analytic step that provides comprehensive spatially resolved maps of tissue heterogeneity (6).

Conventionally, spatial annotation relies on the manual assignment of cell/domain clusters using known marker genes that are readily available from existing studies or databases (7,8). A general workflow begins with the unsupervised clustering of spots based on their transcriptomic profiles; this is followed by an examination of the differentially expressed genes (DEGs) specific to each cluster; and finally, the DEGs are manually matched with known marker genes to assign cell/domain types to spatial spots. This type of workflow requires sufficient knowledge of the biology and markers of the cell/domain types, but it can be time-consuming, labor-intensive and less reproducible (6,9). Moreover, these workflows are sensitive to the choice of clustering methods, presenting challenges in the downstream interpretations (10). An improved strategy for spatial annotation is to automatically annotate the identified clusters using either reference data or leveraging existing information on the cell/domain types. Performing annotations with reference data has been shown to be successful in the context of single-cell RNA sequencing (scRNA-seq) analysis. For example, scmap performs cell annotation by projecting existing reference data with known cell types onto cells in the study data (11). However, the success of this type of analysis relies on the availability of reference data that are ‘similar’ to the study data. On the other hand, the availability of data on cell-type-specific marker genes from existing studies or databases, potentially obtained using either low-throughput or high-throughput systems, further necessitates the efficient utilization of marker-gene information in a ‘qualitative’ manner. To this end, a number of methods have been developed for scRNA-seq data without any consideration of spatial information, including SCINA (12), Garnett (13), CellAssign (14) and scSorter (15). While SCINA and CellAssign use only the expression of marker genes, scSorter and Garnett can utilize information from non-marker genes. Although these methods can be applied to SRT data, they do not consider the invaluable spatial localization information among spots.

To efficiently utilize the existing knowledge based on marker genes for cell/domain types, an ideal annotation method for SRT datasets should be capable of leveraging this ‘qualitative’ information on marker genes with data on non-marker genes while incorporating spatial information to promote spatial smoothness in the cell/domain-type annotation. Because the proportion of non-marker genes is much larger than that of marker genes, non-marker genes also harbor substantial amounts of biological information that can be used to separate cell/domain types. Annotation methods capable

of leveraging marker with non-marker genes can improve our ability to detect spatial cell/domain clusters (14,15). However, the high-dimensional nature of non-marker genes makes the annotation task more challenging and, moreover, requires proper and efficient modeling of this information. Furthermore, for SRT datasets, especially those from tissue sections with laminar structures, for example, brain regions, a desirable spatial annotation method would additionally be able to leverage spatial information.

To address the challenges presented by spatial annotation, we propose the use of a probabilistic model, SpatialAnno, which performs cell/domain-type assignments for SRT data and has the capability of leveraging non-marker genes to assign cell/domain types via a factor model while accounting for spatial information via a Potts model (16,17). To effectively leverage a large number of non-marker genes and overcome the curse of dimensionality, SpatialAnno uniquely models expression levels in a factor model governed by separable cell/domain-type low-dimensional embeddings. As a result, SpatialAnno not only performs spatial cell/domain-type assignments with better accuracy but also estimates cell/domain-type-aware embeddings that can facilitate downstream analyses. We illustrate the benefits of SpatialAnno through extensive simulations and analyses of a diverse range of example datasets collated using different spatial transcriptomics technologies. To show the improved spatial annotation accuracy, we applied SpatialAnno to analyze a 10x Visium datasets for 12 human dorsolateral prefrontal cortex (DLPFC) samples. To illustrate the effectiveness of SpatialAnno in leveraging non-marker genes, we analyzed a mouse olfactory bulb (OB) dataset generated using the ST technology. Using Slide-seqV1/2 datasets for the mouse hippocampus, we demonstrated that SpatialAnno can correctly identify cell-type distribution at near-cell resolution. The utility of SpatialAnno to estimate low-dimensional embeddings is demonstrated by a seqFISH dataset for the mouse embryo.

## Materials and methods

### Model specification

In the SpatialAnno model, we denote  $X$  as the spot-by-gene expression matrix on  $n$  spatial locations. These locations have known spatial coordinates and unknown labels  $y_i, i = 1, \dots, n$ . We can separate genes into a group of  $m$  marker genes and a group of  $p$  non-marker genes, denoted as  $x_{1i} = (x_{i1}, \dots, x_{im})^\top$  and  $x_{2i} = (x_{i,m+1}, \dots, x_{i,m+p})^\top$ , respectively. Suppose prior knowledge of marker genes for  $K$  cell/domain types is encoded as an indicator matrix  $\rho$  of dimension  $m \times K$ , with  $\rho_{jk} = 1$  if gene  $j$  is a maker for cell/domain type  $k$  and 0 otherwise. Following (18–20), we assume that the expression measurements have already been normalized through variance stabilizing transformation and further centered for each gene to have zero mean (see Supplementary Notes).

SpatialAnno models the centered normalized expression vector,  $x_{1i}$ , for marker genes in cell  $i$ , and latent label,  $y_i$ , as

$$\begin{aligned} x_{1i} | y_i = k &\sim \mathcal{N}(\mu_k, \Sigma), \\ \mu_{jk} &= \alpha_j + \rho_{jk}\beta_{jk}, \end{aligned} \quad (1)$$

with the constraint that  $\beta_{jk} \geq 0$ . Here,  $\alpha_j$  is the base expression level for gene  $j$  in the marker group. The intuition is that if gene  $j$  is a marker for cell/domain type  $k$ , then we expect the expression of  $j$  to be higher in these cell/domain types (14) with

an increased magnitude  $\beta_{jk}$ . Note that there is no restriction stating marker genes cannot be expressed in other cell/domain types. We assume the covariance  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_m^2)$ . This simplification significantly reduces the computational cost.

For the high-dimensional non-marker genes, SpatialAnno models their centered normalized expression vector,  $x_{2i}$ , and latent label,  $y_i$ , as

$$\begin{aligned} x_{2i} | z_i &= Lz_i + e_i, \\ z_i | y_i = k &\sim \mathcal{N}(m_k, V), \end{aligned} \quad (2)$$

where factor  $z_i \in R^q$  represents a  $q$ -dimensional embedding of  $x_{2i}$ ;  $L$  is a  $p \times q$  factor loading matrix;  $m_k \in R^q$  is the mean vector for the  $k$ th cell/domain type, and  $V$  is the covariance matrix that is shared across cell/domain types; and  $e_i$  is the residual error and follows an independent normal distribution with mean zero and variance  $\Lambda$ , which is a diagonal matrix, or  $e_i \sim \mathcal{N}(0, \Lambda)$ .

To promote neighborhood similarity in cell/domain types, we follow previous computation (21,22) and assume that cell/domain type  $y_i \in \{1, \dots, K\}$  follows a Potts model characterized by an interaction parameter  $\xi$ , and a neighborhood graph  $\mathcal{S}$ ,

$$p(y | \mathcal{S}, \xi) = \frac{1}{C(\xi)} \exp \left\{ -\xi \sum_{i \sim i'} [1 - I(y_i = y_{i'})] \right\}, \quad (3)$$

where  $i \sim i'$  denotes all neighboring pairs in the neighborhood graph  $\mathcal{S}$ ;  $I(y_i = y_{i'})$  is an indicator function that equals 1 if both the  $i$ th and  $i'$ th locations belong to the same cell/domain type and equals 0 otherwise;  $\xi$  is an unknown interaction parameter that determines the extent of cell/domain type similarity among neighboring locations; and  $C(\xi)$  is the normalizing constant, also known as the partition function that ensures the above probability mass function has a summation of one across all possible configurations of  $y$ .

## Methods for comparison

We compared SpatialAnno with four annotation methods: (i) SCINA (12) implemented in the R package *SCINA* (version 1.2.0), (ii) Garnett (13) implemented in the R package *garnett* (version 0.1.21), (iii) CellAssign (14) implemented in the R package *cellassign* (version 0.99.21) and (v) scSorter (15) implemented in the R package *scSorter* (version 0.0.2). We used the default parameter settings as recommended in their respective tutorials.

SCINA models the expression levels of marker genes using a Gaussian mixture model and enforces a constraint that marker genes should have higher mean expression levels in their corresponding cell types. One key advantage of SCINA is its computational efficiency. CellAssign models the expression count data of marker genes based on a Bayesian probabilistic model that takes into account batch- or sample-specific effects. This approach enhances the accuracy of cell type annotation, particularly when dealing with data from a heterogeneous scRNA-seq population. However, both SCINA and CellAssign rely solely on the expression of marker genes for cell-type annotation. Garnett takes a different approach by first identifying representative cells for known cell types using only marker genes. It then trains a multinomial classifier using all genes with the representative cells and uses this classifier to classify the remaining cells. Garnett also offers a method for rapidly annotating additional datasets by apply-

ing the pre-trained classifier. Notably, non-marker genes are not employed in the initial step of Garnett's approach. In contrast, scSorter combines the expression of both marker genes and non-marker genes. It uses K-means optimization strategies and relies on pre-specified weight parameters to adjust the contribution of marker genes. These weights need to be specified manually.

## Simulations

We performed comprehensive simulations to evaluate the performance of SpatialAnno and compared it with that of alternative annotation methods. The spatial locations of 3639 spots were taken from DLPFC section 151673. Cell/domain types were assigned with manually generated annotations from the original studies (23). We simulated gene expression data for each spot using the *splatter* package (version 1.20.0).

Five marker genes for each cell/domain type were selected from the top DEGs based on log-fold change in expression. We tested the accuracy and robustness of SpatialAnno with the following settings that reflect real-world scenarios.

- I. To test the robustness of SpatialAnno to the erroneous specification of the number of cell/domain types, we considered three scenarios. In the first scenario, marker genes for all seven cell/domain types were provided, and no unknown cell/domain types existed in the expression data. In the second scenario, marker genes for two cell/domain types were removed to create a scenario in which fewer cell/domain types were specified in the marker gene matrix than actually exist in the data. Thus, cells from these two cell/domain types should be assigned to 'unknown'. In the third scenario, the marker genes for nine cell/domain types were added, but two cell/domain types did not appear in the expression data. This mimics a scenario in which more cell/domain types are specified in the marker gene matrix than actually present in the data.
- II. To evaluate the robustness of SpatialAnno to marker gene misspecification, we next created a scenario in which marker genes may be incomplete or incorrect. We randomly flipped a fraction of entries in the binary marker gene matrix  $\rho$  to introduce errors. Specifically, the procedure consisted of two steps. In the first step, a proportion of entries in  $\rho$  that contained one were flipped. In the second step, the same number of entries flipped in the first step was flipped for the entries that contained zero in the original  $\rho$ . The proportions considered were set to 10%, 20% or 30%. Other settings were similar to those in the first scenario in Simulation I.
- III. To assess the capability of SpatialAnno to utilize high-dimensional non-marker genes, we varied the number of non-marker genes as 60, 100, 500, 1000 and 2000. In this setting, we only compared scSorter and Garnett, as only these methods can utilize non-marker genes. Other settings were similar to those in the first scenario in Simulation I.

For each simulation setting, we performed 50 replicate simulations. In each replicate, we applied SpatialAnno and the other methods to annotate each spot.

## Real datasets

### Human dorsolateral prefrontal cortex data generated using 10x Visium

We downloaded a human DLPFC dataset (23) generated using the 10x Visium platform from <http://spatial.libd.org/spatialLIBD/>. In this dataset, there were 12 tissue sections, which contained a total of 33 538 genes measured on average over 3973 spots. We used the sample ID151673, which contains expression measurements of 33 538 genes on 3639 spots, as the main analysis example. We presented the results for the other 11 samples in the Supplementary Figures. For all the sections, we extracted the top 2000 spatially variable genes with SPARK-X (24) before performing annotations. To identify layer-specific marker genes for annotation, we used tissue section 15 1507 as the reference data. This dataset contained 33 538 genes for 4226 spots. For each layer, the top 5 DEGs were selected as its marker genes. The final marker gene list is available in Supplementary Table S1.

### Mouse olfactory bulb data by spatial transcriptomics (ST)

We obtained the mouse olfactory bulb ST data from the spatial transcriptomics research website (<https://www.spatialresearch.org/>). These data consist of gene expression levels in the form of read counts that were collected for a number of spatial locations. We followed the methods of previous studies (25,26) to focus on the mouse OB Section 12, which contains 16 034 genes and 282 spatial locations. We presented the results for the other 11 sections in the Supplementary Figures. We extracted the top 3000 most highly variable genes with function SCTransform implemented in *Seurat* (version 4.0.5) (27) before performing annotations. To construct the marker gene list for annotation, we perform differential expression analysis on scRNA-seq data (28) from the Gene Expression Omnibus (GEO; accession number GSE121891). This scRNA-seq data was collected from the mouse olfactory bulb and contains 18 560 genes and 12 801 cells for five cell types: granule cells (GC,  $n = 8614$ ), olfactory sensory neurons (OSNs,  $n = 1200$ ), periglomerular cells (PGC,  $n = 1693$ ), mitral and tufted cells (M-TC,  $n = 1133$ ), and external plexiform layer interneurons (EPL-IN,  $n = 161$ ). For each cell type, the top four DEGs were selected as its marker genes. The final marker gene list is available in Supplementary Table S2.

### Mouse hippocampus Slide-seq data and Slide-seqV2 data

We obtained the mouse hippocampus Slide-seq dataset and Slide-seqV2 dataset (3) from the Broad Institute's Single Cell Portal ([https://singlecell.broadinstitute.org/single\\_cell/study/SCP948/robust-decomposition-of-cell-type-mixtures-in-spatial-transcriptomics](https://singlecell.broadinstitute.org/single_cell/study/SCP948/robust-decomposition-of-cell-type-mixtures-in-spatial-transcriptomics)). The Slide-seq dataset consists of gene expression measurements in the form of read counts for 22 457 genes and 34 199 spatial locations. The Slide-seqV2 dataset consists of gene expression measurements in the form of read counts for 23 264 genes and 53 208 spatial locations. In the analysis, we filtered out genes that had fewer than 20 counts on all locations and filtered out locations that had fewer than 20 genes with nonzero counts. These filtering criteria led to final sets of 14 481 genes and 31 664 cells for Slide-seq dataset, and 16 121 genes and 51 212 cells for Slide-seqV2 dataset. In addition, for both datasets, we extracted the top 2000 most spatially variable genes with SPARK-X (24) before performing annotations. To construct marker genes for annotation, we obtained the DropViZ scRNA-seq dataset

(29) from the Broad Institute's Single Cell Portal. This data was collected from the mouse hippocampus, which contained 22 245 genes and 52 846 cells for 19 cell types. For each cell type, the top five DEGs were selected as marker genes. Besides the 19 cell types, we added another two cell types, *Slc17a6* neurons and Hb neurons, and their marker genes were extracted from the original study (29). The final marker gene list used is available in Supplementary Table S3.

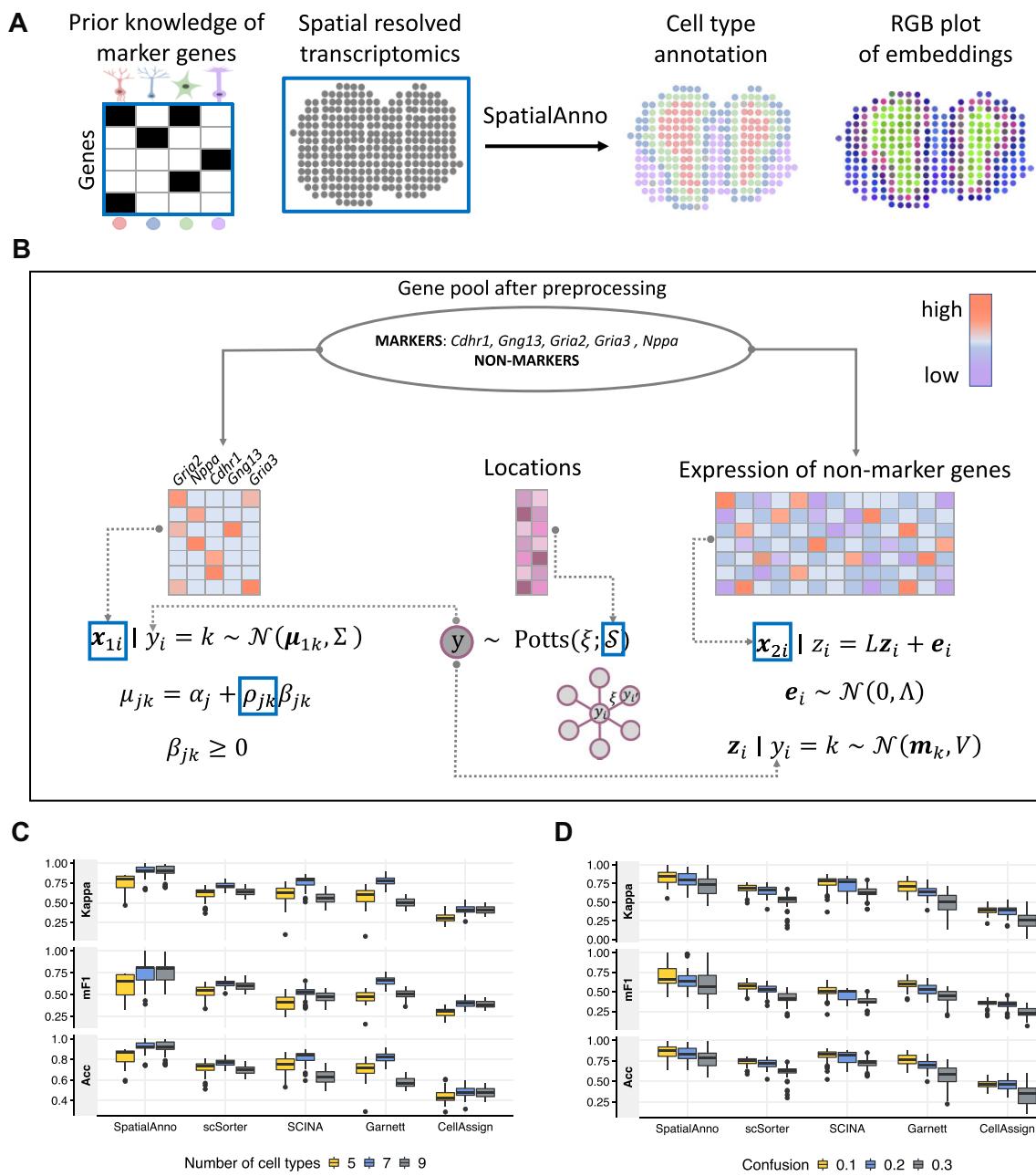
### Mouse embryo data by seqFISH

We obtained the mouse embryo seqFISH data (2) from <https://marionilab.cruk.cam.ac.uk/SpatialMouseAtlas/>. This dataset profiles the expression of 387 selected target genes from three mouse embryo tissue sections. Cell segmentation is performed using a combination of aligning membrane stains to the first hybridization round and training them with a machine learning toolkit called ilastik (30). This process generates probability maps, which are then used to create 2D-labeled cells for each z slice. mRNA transcript signals are located by finding local maxima above a threshold, and these spots are assigned to corresponding cells based on location, generating a gene-cell count matrix. In total, the dataset includes gene-cell count matrices for 19 451, 14 891 and 23 194 cells, respectively. We calculated normalized expression log counts for each cell using `logNormCounts` function in the R package *scuttle* (31) with cell-specific size factors. To construct the marker gene list, we used Embryo 3 as a reference; this dataset contains 24 cell types. For each cell type, the top eight DEGs were selected as marker genes. We removed marker genes for two cell types, ExE endoderm cells and blood progenitors, as there were too few (<30) of these cells. The cell type 'Low quality' was also removed. The final marker gene list used contained 21 cell types and is available in Supplementary Table S4.

### Evaluation metrics

We evaluated annotation performance using three metrics, that is, Kappa, mF1 score and ACC, as suggested in previous single-cell data annotation studies (11,32). ACC was defined as the proportion of spots that were classified into the correct types. Kappa is generally thought to be a more robust measure than ACC, since it takes into account the possibility that the agreement occurs by chance. The cell-level F1 score considers each cell to be an individual classification task with a true cell-type assignment (and potentially multiple incorrect cell-type assignments) for the purposes of calculating precision and recall (Supplementary Notes).

We also compared the low-dimensional embeddings estimated in SpatialAnno with those from PCA and DR-SC (22). In detail, we first extracted the top 15-dimensional components and then summarized those top components as three tSNE components and visualized the resulting tSNE components with RGB colors in the RGB plot. To show that the estimated embeddings carry the most information about cell/domain types, we evaluated the conditional correlation coefficients between the true cell/domain labels and the observed gene expression, given the estimated embeddings in SpatialAnno. Furthermore, the embeddings in SpatialAnno improve clustering performance. With embeddings from SpatialAnno, PCA and DR-SC, we performed clustering analysis using the Louvain community detection algorithm imple-



**Figure 1.** Schematic overview of SpatialAnno and its performance in simulation studies. **(A)** SpatialAnno employs spatial transcriptomics data along with a known marker-gene list in its analysis. With these two datasets as input, SpatialAnno performs spatial annotation via a probabilistic model that combines both marker and non-marker gene expression data, and produces both domain/cell-type assignments and low-dimensional embeddings for all spatial locations as output. **(B)** Overview of the SpatialAnno probabilistic model. Latent cell/domain types (shown in the gray circle) and observed data (shown in the blue boxes) are shown along with the distributional assumptions. **(C)** Kappa, mF1, and ACC of SpatialAnno, scSorter, SCINA, Garnett and CellAssign for simulation data from seven cortical layers; different numbers of cell/domain types are provided as a list of marker genes. **(D)** Kappa, mF1 and ACC of SpatialAnno, scSorter, SCINA, CellAssign and Garnett for simulation data from seven cortical layers; different proportions of marker genes are erroneously specified.

mented in the R package *Seurat* (version 4.1.1), and evaluated clustering performance using the ARI (33).

## Results

### Overview of SpatialAnno

Similarly to other methods that assign known cell/domain types to cells using information about marker genes, SpatialAnno takes as input normalized gene expression matrix,

spatial location information, and a list of marker genes for known cell/domain types (Figure 1A). SpatialAnno automatically performs cell/domain-type assignments while providing low-dimensional embeddings for all spatial spots. Although most *in situ* capturing technologies have limited spatial resolution, with each measured location possibly containing multiple cell types, SpatialAnno can still provide a crucial understanding of tissue organization by annotating domain types. When marker genes for known cell types are available, Spa-

tialAnno can be used to annotate cell types for measured locations. However, it is important to note that relying solely on the major cell type identified in each location could potentially lead to biased results. Based on the latent cell/domain type for each spot, SpatialAnno builds a ‘semi-supervised’ Gaussian mixture model to modulate the over-expression of marker genes and a hierarchical factor model to relate non-marker gene expression to the cell/domain separable latent embeddings while accounting for the spatial smoothness of the cell/domain types with a Potts model (Figure 1B). Uniquely, SpatialAnno, via the factor model, allows for the assignment of cell/domain types that leverage a large number of non-marker genes, and, via the Potts model, is more likely to assign the same cell/domain type to neighboring spots, promoting spatial smoothness in the cell/domain types. Notably, with expression data for both marker and non-marker genes, SpatialAnno simultaneously assigns each spot known cell/domain types while obtaining low-dimensional embeddings for each spot, which can facilitate other downstream analyses. Similarly to other methods, SpatialAnno automatically labels spatial spots that do not belong to any known cell/domain type as ‘unknown’, preventing incorrect assignment when novel cell/domain types are present.

### Validation using simulated data

We conducted simulations to evaluate the performance of SpatialAnno and compared the results with those of non-spatial annotation methods commonly applied to scRNA-seq data: SCINA, Garnett, CellAssign, and scSorter. Briefly, we simulated gene expression counts using a splatter model (34) for seven cortical layers using labels from the DLPFC data. Then, we selected five marker genes for each layer based on the log-fold change in expression (see Supplementary Notes). In total, we obtained 35 marker genes and 2000 non-marker genes for 3639 spots from seven layers. For each simulated SRT dataset, we applied SpatialAnno and the four other methods to perform spatial domain annotation. We used Cohen’s Kappa, mean F1 (mF1) score, and classification accuracy (ACC) (see Supplementary Notes) to quantify the concordance between the detected spatial domains and the seven labeled cortical layers (11,14). We performed 50 replicate simulations for each setting. To determine if the three performance measures of the compared methods were distinguishable, we computed the Bayes factor (35,36) to directly compare the performance of each method against SpatialAnno. A Bayes factor  $>3$  was considered statistically different. (36).

When the correct number of layers was specified, SpatialAnno ( $\text{Kappa} = 0.903$ ,  $\text{mF1} = 0.807$  and  $\text{ACC} = 0.922$ ) outperformed all other methods in terms of annotation accuracy (Figure 1C; number of cell/domain types = 7), with Bayes factors  $>10$  (Supplementary Figure S1A). After varying the number of cell/domain types with marker genes, the SpatialAnno annotation still outperformed all other methods (Figure 1C; number of cell/domain types = 5 or 9), with Bayes factors  $>10$  (Supplementary Figure S1A). Unsurprisingly, SpatialAnno performed worse when there were five cell/domain types with marker genes ( $\text{Kappa} = 0.839$ ,  $\text{mF1} = 0.729$  and  $\text{ACC} = 0.883$ ) than seven or nine ( $\text{Kappa} = 0.900$ ,  $\text{mF1} = 0.803$  and  $\text{ACC} = 0.918$ ). The latter two cases (seven and nine cell/domain types) led to comparable annotation performances for SpatialAnno and CellAssign. In contrast, annotation performance decreased for the other methods when we

included marker genes for irrelevant cell/domain types. We examined the robustness of SpatialAnno when there were various degrees of marker gene misspecification (Figure 1D), as well as the presence of shared marker genes across different cell types. As the proportion of misspecified marker genes increased, the annotation performance decreased for all methods, but SpatialAnno still outperformed all other methods in terms of annotation accuracy (Kappa, mF1 and ACC), with Bayes factors  $>3$  (Supplementary Figure S1B). Similarly, when the proportion of shared marker genes across different cell types increased, we observed consistent outcomes (Supplementary Figure S1C and S1D).

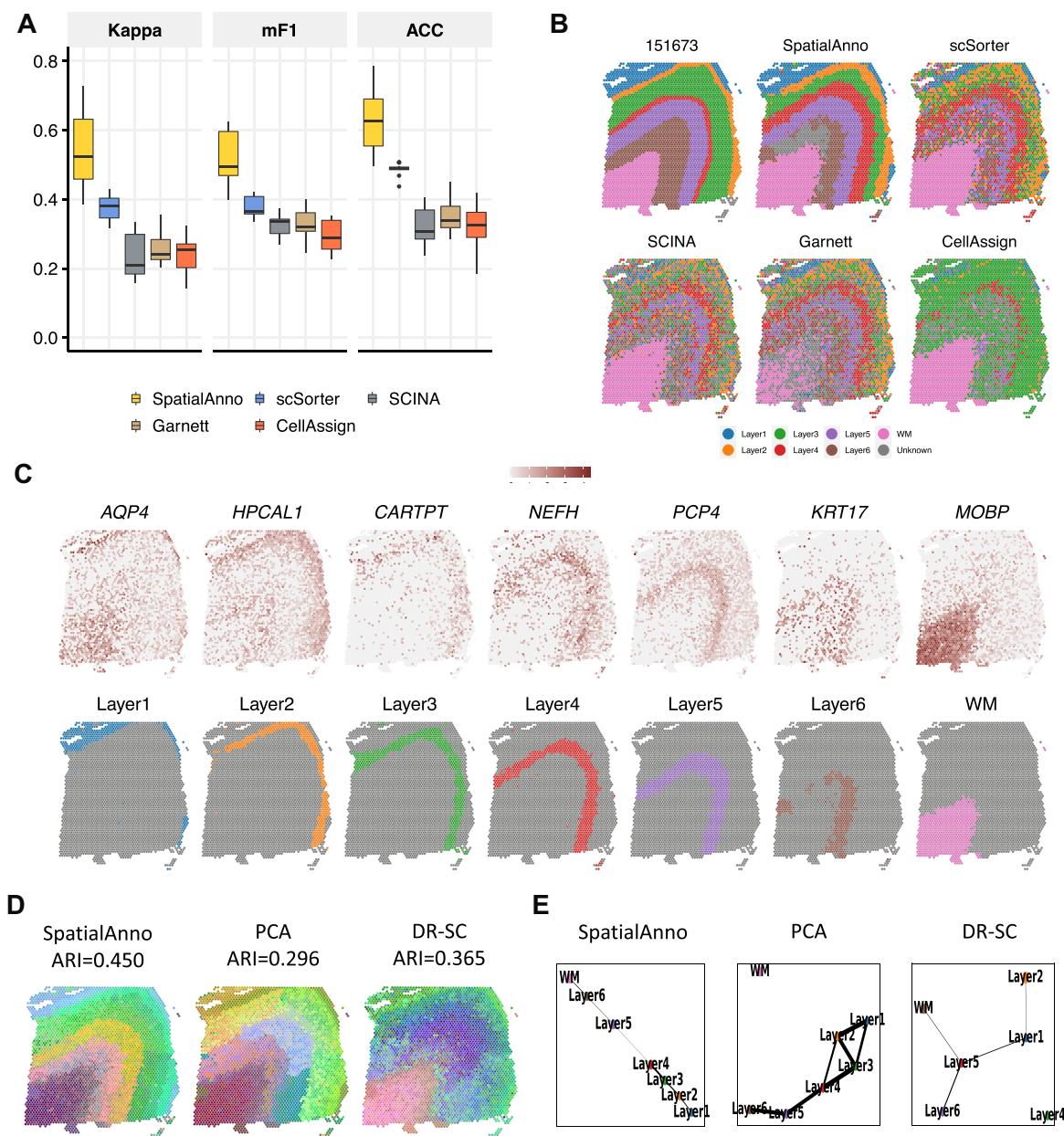
Next, we examined the effectiveness of SpatialAnno, which leverages various amounts of non-marker information compared with the scSorter and Garnett methods, also capable of leveraging non-marker genes (Supplementary Figure S2A). As the number of non-marker genes increased from 60 to 2000, SpatialAnno showed 10.3%, 21.9% and 8.1% improvements in annotation accuracy for Kappa, mF1 and ACC, respectively, while the annotation accuracies of scSorter and Garnett were almost unchanged, with the changes being -0.6% and -0.6% for Kappa, 1.7% and -0.1% for mF1, and -0.1% and -0.7% for ACC, respectively. These results suggest that SpatialAnno can effectively leverage various numbers of non-marker genes.

In addition to the spatial spots being accurately annotated, the low-dimensional embedding of non-marker genes from SpatialAnno was cell/domain-type informative. Clustering performance using low-dimensional embeddings with either marker genes or non-marker genes, or a combination of the two, with a comparable adjusted rand index (ARI) between marker and non-marker genes, is shown in Supplementary Figure S2B and C. Not surprisingly, combining both embeddings for marker and non-marker genes led to improved ARIs in all scenarios, demonstrating the benefits of borrowing information from non-marker genes when annotating cell/domain types. In addition, the Pearson’s correlation coefficients for the relationship between the observed expression and the estimated labels, given the embeddings from SpatialAnno, were much smaller than those for the principal component analysis (PCA), but comparable to those for the DR-SC (22) (Supplementary Figure S2D and E). These results suggest SpatialAnno embeddings can capture cell/domain-type-relevant information for each spot, thus facilitating the downstream analysis.

Finally, we evaluated the computational efficiency of all methods for different numbers of cell/domain types, as shown in Supplementary Figure S2F. SpatialAnno was computationally efficient and comparable in efficiency to SCINA and scSorter, and all three were faster than Garnett and CellAssign.

### SpatialAnno improves annotations of known layers in human dorsolateral prefrontal cortex

We applied SpatialAnno and the four methods to the analysis of human dorsolateral prefrontal cortex (DLPFC) 10x Visium data (23). In this dataset, there were 12 tissue sections from three adult donors with a median depth of 291 million reads for each sample, a median of 3844 spatial spots per section and a mean of 33 538 genes per spot (Supplementary Table S5). Based on a manual examination of cytoarchitecture and specific marker genes, each tissue section was carefully annotated by the original study (23) in one of the six layers of the prefrontal cortex or white matter (WM). Taking sample



**Figure 2.** Spatial domain annotation in the DLPFC 10x Visium dataset. **(A)** Boxplots of Kappa, mF1, and ACC showing the accuracy of different methods for domain annotation across 12 tissue sections. **(B)** Spatial domain annotation in tissue sample ID151673 for ground truth, SpatialAnno, scSorter, SCINA, Garnett and CellAssign. **(C)** Top, expression levels of corresponding layer-specific marker genes. Bottom, annotations by SpatialAnno are shown on each spot. **(D)** RGB plots for low-dimensional embedding inferred by SpatialAnno, PCA and DR-SC. As end-to-end annotation approaches, scSorter, SCINA, Garnett and CellAssign cannot be utilized to extract low-dimensional embeddings. **(E)** PAGA graphs generated by SpatialAnno, PCA and DR-SC embeddings for DLPFC Section ID151673.

ID151507 as a reference, we constructed a marker-gene list that contained five marker genes for each of the seven layers (see Supplementary Notes).

Taking manual annotations as ground truth, we first evaluated the performance of spatial annotation using Kappa, mF1, and ACC for each of the 12 tissue sections (Figure 2A). SpatialAnno annotated spatial domains more accurately (median Kappa = 0.524, median mF1 = 0.494 and median ACC = 0.628) than scSorter (median Kappa = 0.381, median mF1 = 0.366 and median ACC = 0.489), SCINA (median Kappa = 0.209, median mF1 = 0.337 and median ACC = 0.307), Garnett (median Kappa = 0.24, median mF1 = 0.32 and me-

dian ACC = 0.339) and CellAssign (median Kappa = 0.253, median mF1 = 0.29 and median ACC = 0.326), with Bayes factors >50 (Supplementary Figure S3). The heatmap of the spatial assignments from SpatialAnno and the other methods and the manual annotations for sample ID151673 are shown in Figure 2B. SpatialAnno achieved the best annotation accuracy (Kappa = 0.634, mF1 = 0.619 and ACC = 0.685), while the annotations from scSorter, SCINA and CellAssign were only accurate for the WM, and Garnett completely failed to assign the WM region. Notably, the domains identified in SpatialAnno were spatially smooth, continuous, and well matched with the elevated expression levels of marker genes for each

layer (Figure 2C and Supplementary Figure S4–S15), such as *PCP4* and *MOBP* that are marker genes for layer 5 and WM, respectively (23,37).

To evaluate the robustness of SpatialAnno, we obtained marker genes from the other DLPFC tissue section that contained seven layers and performed spatial annotation for the remainder of the 11 tissue sections (see Supplementary Notes). Using the top 5/10/15 DEGs as marker genes for each layer, SpatialAnno achieved the best annotation accuracy according to Kappa, mF1 and ACC. The annotation accuracies of all other methods for the other tissue sections were slight worse than for those when sample ID151507 was used as a reference (Supplementary Figure S16A), which is consistent with the simulations involving the misspecification of marker genes (Figure 1D). This suggests that annotation accuracy can be impaired when inaccurate marker genes are used. However, this difference became negligible when the number of marker genes for each layer was 15. Furthermore, we examined the robustness of SpatialAnno using marker genes for irrelevant cell types, those not present in the studied SRT dataset. For samples ID151669–151672 from Donor 2, which only contained five cortical layers, we applied SpatialAnno and other methods using marker genes for the seven layers. As shown in Supplementary Figure S16B, SpatialAnno achieved the best annotation performance for these samples.

Uniquely amongst the methods, SpatialAnno's estimated embeddings were highly informative for the DLPFC layers in the 12 sections. The clustering accuracies, determined using the ARI for embeddings from marker, non-marker, and a combination of the two, respectively, were shown in Supplementary Figure S16C, with the largest ARI value for embeddings from a combination of the two. Clearly, embeddings from non-marker genes harbored substantial amount of information about spatial domains, even more than the marker genes. When using a combination of marker and non-marker genes, the embeddings led to improved clustering performance, suggesting that annotation based on both marker and non-marker genes improved the annotation accuracy. Red/green/blue (RGB) plots using three tSNE components for the embeddings in sample ID151673 estimated by SpatialAnno revealed a more clear laminar structure for DLPFC than those by PCA or DR-SC (Figure 2D). Such stronger structure predictivity from SpatialAnno is numerically supported by its higher ARI (0.450) compared to PCA (ARI = 0.296) and DR-SC (ARI = 0.365). Moreover, an estimated PAGA graph (38) using SpatialAnno embeddings demonstrated the almost linear development trajectory from WM to layer 1, while the PAGA graphs using both PCA and DR-SC embeddings were less clearly delineated (Figure 2E and Supplementary Figure S4–S15). To better understand the impact of each component in SpatialAnno, we conducted additional experiments on the DLPFC dataset. As shown in Supplementary Figure S16D, we demonstrate the performance of the model when one or more components were disabled.

### SpatialAnno correctly identifies cells in mouse olfactory bulb

To quantitatively demonstrate the performance of SpatialAnno compared with SCINA, scSorter, CellAssign and Garnett in domain-type annotation, we analyzed a mouse OB data generated using ST technology. This dataset represented 12 tissue sections with a median of 16 024 gene expression

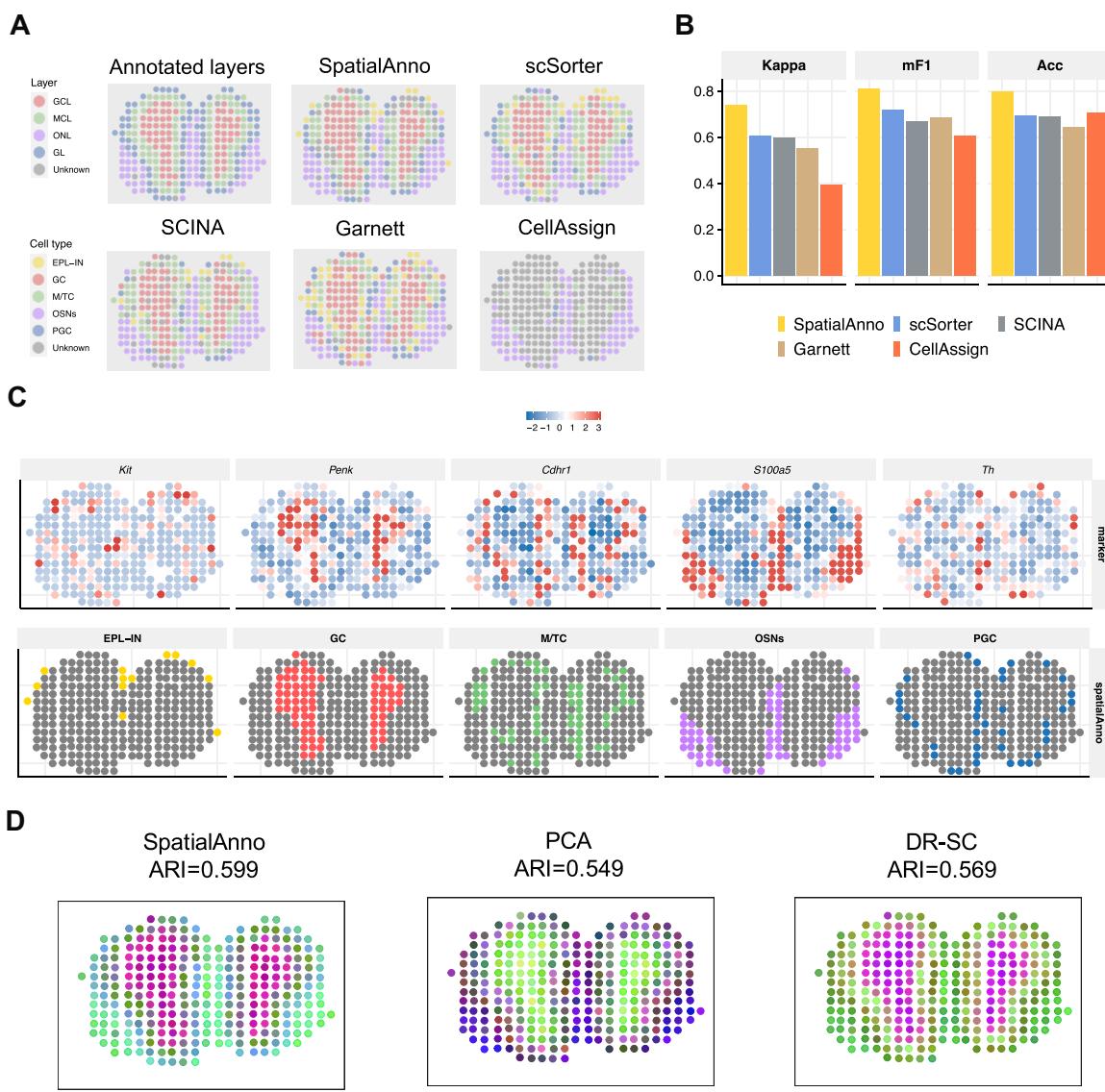
measurements among a median of 266 spots (Supplementary Table S6).

Taking the four anatomic layers manually annotated based on H&E staining as ground truth (Figure 3A), we first evaluated the performance of the spatial annotation using Kappa, mF1 and ACC for Section 12 (Figure 3B). SpatialAnno annotated spatial domains more accurately (Kappa = 0.739, mF1 = 0.812 and ACC = 0.800) than scSorter (Kappa = 0.608, mF1 = 0.718 and ACC = 0.696), SCINA (Kappa = 0.598, mF1 = 0.670 and ACC = 0.689), CellAssign (Kappa = 0.395, mF1 = 0.607 and ACC = 0.707) and Garnett (Kappa = 0.552, mF1 = 0.686 and ACC = 0.646). We examined the robustness of SpatialAnno by including marker genes for two irrelevant cell types (endothelial and mural cells) that were not present in this section, and SpatialAnno achieved the best annotation performance (Supplementary Figure S17A). To illustrate the effectiveness of leveraging non-marker information, we evaluated the performance of the spatial annotation by SpatialAnno, scSorter, and Garnett with 30, 300 or 3000 non-marker genes, as only these three methods are able to leverage non-marker gene information. SpatialAnno achieved higher annotation accuracy when more non-marker genes were used, while the difference in performance between 300 and 3000 non-marker genes was minimal for SpatialAnno (Supplementary Figure S17B). In contrast, scSorter and Garnett performed similarly with 30 or 300 non-marker genes, but their performance deteriorated when 3000 non-marker genes were applied.

SpatialAnno recovered the laminar structure of the mouse OB across 12 sections (Supplementary Figure S18). The mouse OB has a multi-layered cellular architecture in the order, from the inner to outer layer, of granule cell layer (GCL), mitral cell layer (MCL), glomerular layer (GL) and the nerve layer (ONL). Detailed assignments by SpatialAnno and the other four methods for Section 12 are shown in Figure 3A. The cell types annotated by SpatialAnno accurately represented this laminar structure, while CellAssign incorrectly assigned ‘unknown’ cells to regions belonging to GCL, MCL and GL. Moreover, the annotation patterns of Garnett were rather chaotic, while scSorter and SCINA failed to distinguish periglomerular cells (PGC) in the GL.

We further examined the expressions of marker genes specific to each layer, including *Kit* for external plexiform layer interneuron (EPL-IN) (28), *Penk* for granule cells (GC) (39), *Cdhr1* for mitral and tufted cells (M/TC) (40), *S100a5* for olfactory sensory neurons (OSN) (41) and *Th* for PGC (42) (Figure 3C). Although the three methods provided similar assignments for GC, M/TC, OSN and PGC, their assignments for EPL-IN were quite different. EPL-IN are located adjacent to GL in the external plexiform layer comprises PGC (28). SpatialAnno assigned spots near PGC to EPL-IN; however, scSorter and Garnett did not (Supplementary Figure S19). As the ground truth for the EPL-IN locations was unknown, we manually combined the inferred EPL-IN with the adjacent layers in different ways: (i) by combining the inferred EPL-IN and PGC and (ii) by combining the inferred EPL-IN, M/TC and PGC. SpatialAnno still achieved the best annotation accuracy (Supplementary Figure S17C & D).

Another key benefit of SpatialAnno is its ability to extract low-dimensional embeddings relevant to different cell types from the high-dimensional non-marker genes, which is useful for many downstream analyses. We summarized the low-dimensional embeddings inferred by SpatialAnno (Sup-



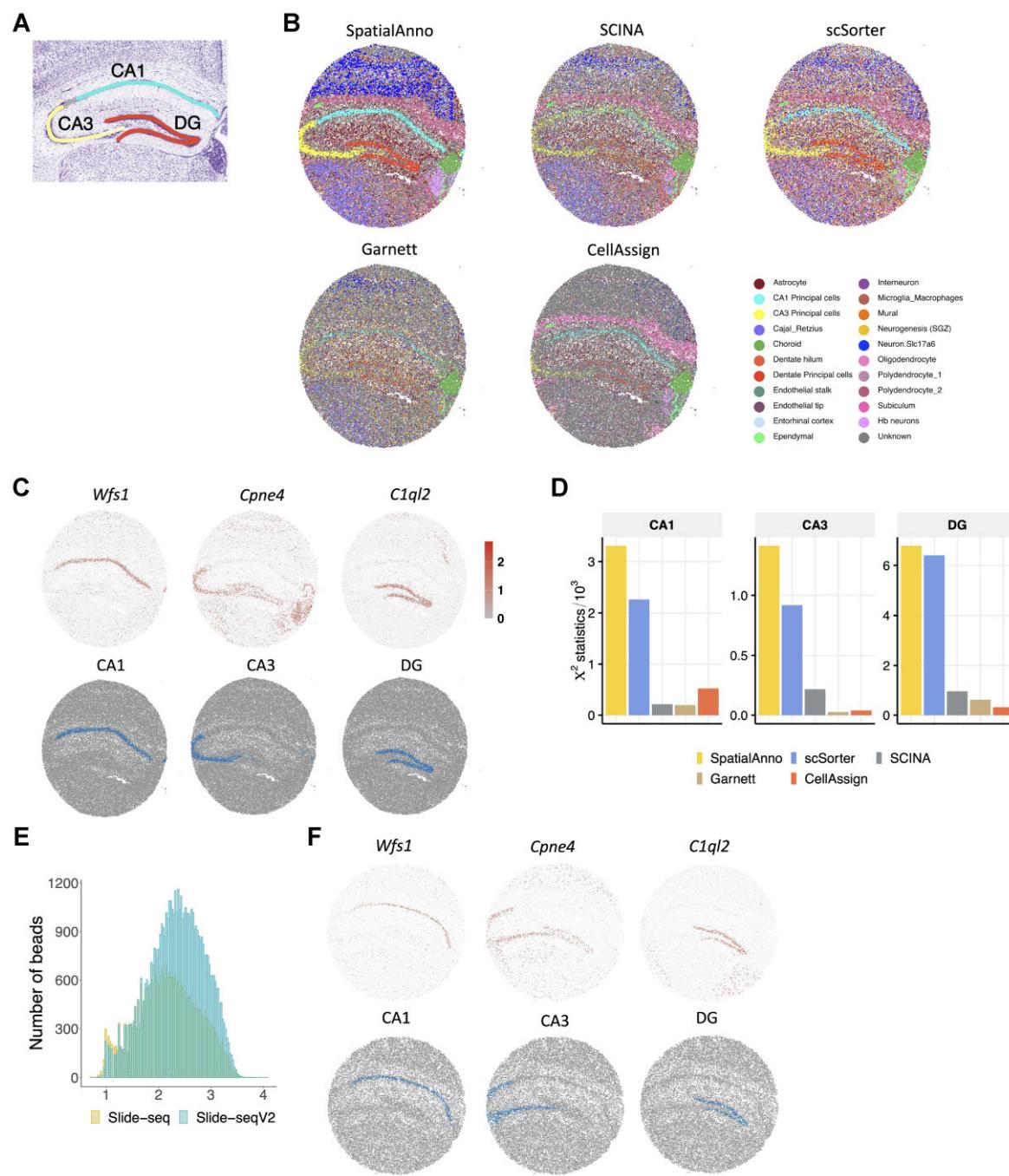
**Figure 3.** Spatial annotation in the mouse olfactory bulb dataset. **(A)** Anatomic layers annotated based on H&E staining of the olfactory bulb, and cell-types inferred by SpatialAnno, scSorter, SCINA, Garnett and CellAssign. **(B)** Bar plots of Kappa, mF1 and ACC showing the domain-type annotation accuracy of different methods. **(C)** Top, expression levels of corresponding cell-type-specific marker genes. Bottom, annotations by SpatialAnno are shown on each spot. **(D)** RGB plots of low-dimensional embeddings inferred by SpatialAnno, PCA and DR-SC. As end-to-end annotation approaches, scSorter, SCINA, Garnett and CellAssign cannot be utilized to extract low-dimensional embeddings.

plementary Figure S17E), PCA and DR-SC into 3D tSNE components and visualized the resulting components in the RGB plot. The RGB plot (Figure 3D) shows the multi-layered architecture of the mouse OB, with neighboring spots sharing more similar colors to those farther away. To compare the predictive powers of these low-dimensional embeddings for the four anatomic layers annotated based on H&E staining, we applied the Louvain community detection algorithm to spot clustering using the *Seurat* R package. The clusters identified by SpatialAnno depicted the multi-layered structures more accurately (ARI = 0.599) than those of PCA (ARI = 0.549) or DR-SC (ARI = 0.569).

**SpatialAnno reveals cell-type distribution in mouse hippocampus with SRT data at near-cell resolution**  
To show the cell-type distribution in the mouse hippocampus, we applied SpatialAnno and the other methods to

the analysis of a mouse hippocampus dataset generated using Slide-seqV2, which quantifies transcriptome-wide expression levels at near-cellular resolution with 10- $\mu$ m barcoded beads (3). This dataset contains expressions for 23 264 genes over 53 208 spatial locations (Supplementary Table S7). As shown in the Allen Reference Atlas (Figure 4A), the primary regions in the mouse hippocampus were composed of the cornu ammonis (CA1-3) and dentate gyrus (DG).

SpatialAnno clearly identified a ‘cord-like’ structure as well as an ‘arrow-like’ structure in the hippocampal subfields in CA1, CA3 and DG (Figure 4B), which is consistent with the annotation of hippocampus structures in the Allen Reference Atlas (Figure 4A). In contrast to SpatialAnno, the other methods SCINA, Garnett, and CellAssign showed blurred/incorrect localizations for the primary hippocampal subfields in CA3 and DG and were unable to reveal the main



**Figure 4.** Spatial cell-type annotation of the mouse hippocampus dataset. **(A)** Annotation of hippocampus structures from the Allen Reference Atlas of an adult mouse brain. **(B)** Spatial annotation of the Slide-seqV2 hippocampus section by SpatialAnno, scSorter, SCINA, Garnett and CellAssign. **(C)** Top, expression levels of corresponding cell-type-specific marker genes. Bottom, annotations by SpatialAnno of the Slide-seqV2 hippocampus section are shown on each spot. The examined cell types were CA1 cells, CA3 cells and dentate cells. **(D)** Results of Pearson's chi-squared test of correlation between expression patterns of marker genes and the three hippocampal subfields identified by different methods. **(E)** Total UMIs per bead for Slide-seq (yellow,  $n = 34$ , 199 spots) versus Slide-seqV2 (blue,  $n = 53$ , 208 spots) in the mouse hippocampus sections. **(F)** Top, expression levels of corresponding cell type specific marker genes. Bottom, annotation by SpatialAnno of the Slide-seq hippocampus section is shown on each spot.

structures of the mouse hippocampus (Figure 4B and Supplementary Figure S20–S22). The hippocampal subfields identified by scSorter were surrounded by a blurry border, with many different cell types allocated to the same region. Additionally, all the methods except SpatialAnno failed to accurately allocate the habenula (Hb) neurons, which should reside left to and below the choroid plexus. Careful examination of marker genes further demonstrated the superior accuracy

of SpatialAnno (Figure 4C), i.e., *Wfs1*, *Cpne4* and *C1ql2* for CA1, CA3 and DG, respectively.

We quantified the annotation performance of the different methods by examining the correlations between the expression patterns of the marker genes and the three hippocampal subfields identified by the different methods. Pearson's chi-squared test demonstrated a substantial improvement in the magnitude of associations provided by SpatialAnno (Fig-

ure 4D). The RGB plot for SpatialAnno displayed clear regional segregation of the hippocampus (Supplementary Figure S23A). Specifically, compared with the RGB plots for PCA and DR-SC, the plot for SpatialAnno clearly depicted the Hb region.

Finally, we validated the cell-type distributions identified for an independent slide from the mouse hippocampus profiled using Slide-seq. As with the initial version of Slide-seqV2, the transcript detection sensitivity of Slide-seq is relatively low (Figure 4E). SpatialAnno successfully identified the hippocampal subfields in this Slide-seq data (Supplementary Figure S23B–D and Supplementary Figure S24–S26). The annotated regions for CA1, CA3 and DG with their marker gene expressions are shown in Figure 4F.

### Embeddings estimated by SpatialAnno lead to biologically relevant trajectories in mouse embryo

We further applied SpatialAnno and the other methods to the analysis of a dataset obtained from three mouse embryo sections collated at the 8–12 somite stage using seqFISH (2), which has the capability of probing the expression of a targeted gene set at the single-cell resolution by image processing and single-cell segmentation (2). Each of the three mouse embryo sections contained expression level measurements for 351 genes, chosen to recover the cell-type identities at these developmental stages, from around 20 000 cells, as well as their physical locations (Supplementary Table S8). After selecting 168 marker genes for 21 cell types (see Supplementary Notes), 183 non-marker genes remained for annotation analysis.

The original study provided manual annotations for the cells based on their nearest neighbors in the Gastrulation atlas (43). For each method, we summarized the annotation accuracy using both Kappa, mF1 and ACC for each embryo section (Figure 5A and Supplementary Figure S27). SpatialAnno achieved the highest Kappa, mF1 and ACC in two of the three sections and was only surpassed by CellAssign for the second embryo section. For Embryo 1, the annotations of different methods are shown in Figure 5B. Clearly, cell-type distributions identified by SpatialAnno were well matched with the expression of their corresponding marker genes (Figure 5C).

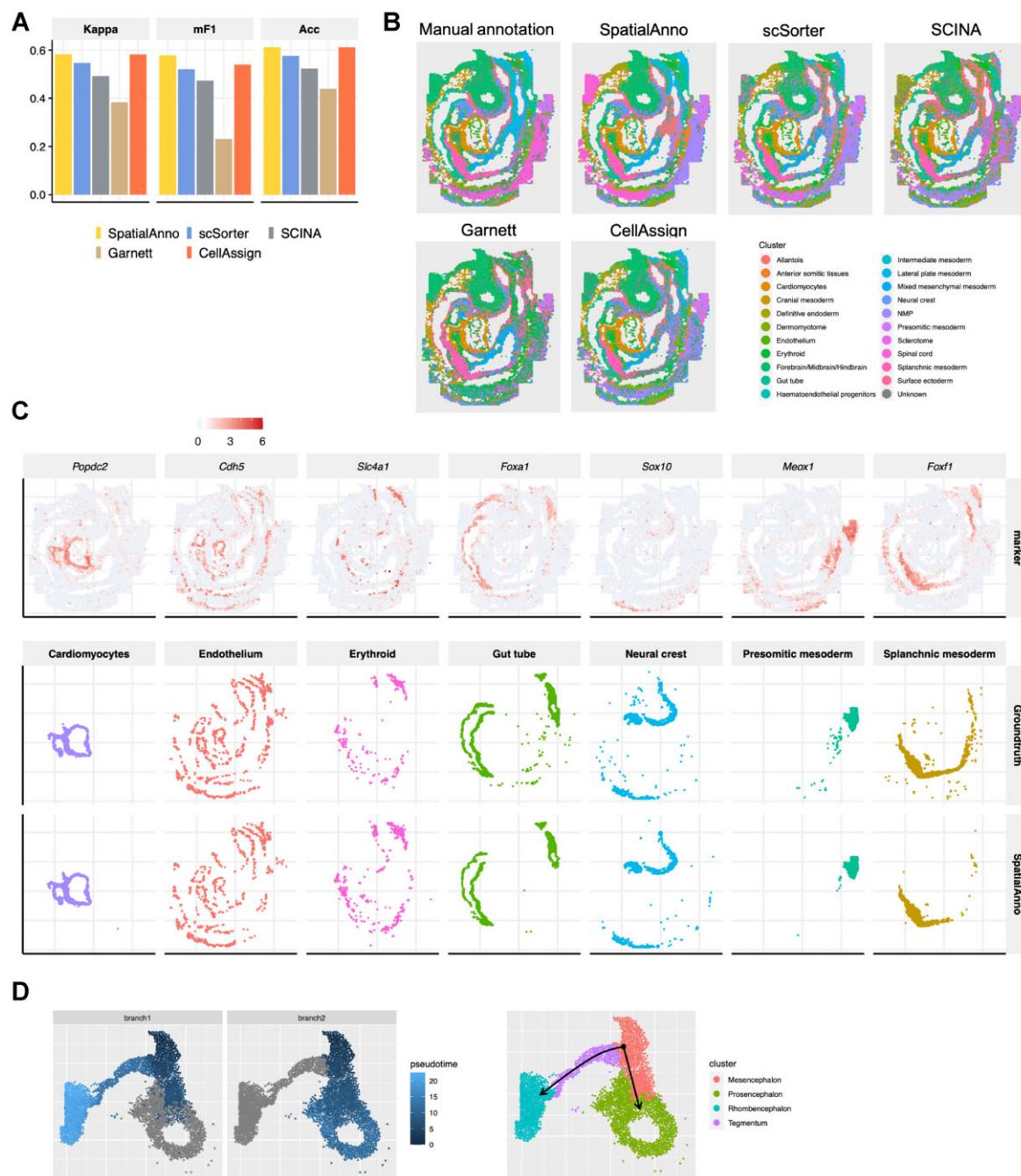
For the embeddings uniquely estimated by SpatialAnno, we performed trajectory inference on brain cells to investigate the spatiotemporal development of the mouse brain and detected two linear trajectories (Figure 5D). We observed the lowest pseudotime values in the mesencephalon, which diffused smoothly toward the tegmentum followed by the rhombencephalon in one branch, and towards the prosencephalon in another branch (Figure 5D). More importantly, the diffusion patterns were spatially continuous and smooth. The detected trajectories delineated the spatial trajectories of mouse brain development, which are in agreement with the findings of recent studies (2). In contrast, the trajectories identified using embeddings from either PCA or DR-SC lacked spatial continuity (Supplementary Figure S28A and B). We further examined genes associated with the inferred pseudotime, and a heatmap of the expression levels of the top 20 significant genes suggested that there were interesting expression patterns over pseudotime (Supplementary Figure S28C). A mesencephalon and prosencephalon marker gene, *Otx2* (44,45), showed higher expression levels in the early stage of development, while at a later stage, its expression levels were sub-

stantially suppressed (Supplementary Figure S28D). In contrast, the expression levels of a gene enriched in the rhombencephalon, *Sfrp1* (46), changed from low to high (Supplementary Figure S28D). These results concur with the formation of the midbrain-hindbrain boundary (47,48), and this is supported by the observation that these two genes could be used to identify the precise boundary between the mesencephalon and rhombencephalon (Supplementary Figure S28E).

## Discussion

SpatialAnno takes, as input, the normalized gene expression matrix, the physical location of each spot, and a list of marker genes for known cell/domain types. The output of SpatialAnno comprises the estimated posterior probability of each spot belonging to each cell/domain type and the low-dimensional embeddings of each spot for non-marker genes. To efficiently capitalize on both marker and non-marker genes, SpatialAnno uniquely models the expression levels of non-marker genes via a factor model governed by cell/domain-type separable low-dimensional embeddings and simultaneously promotes spatial smoothness via a Potts model. As a result, SpatialAnno provides improved spatial cell/domain-type assignments, and its estimated low-dimensional embeddings are cell-type-relevant and can facilitate downstream analyses such as trajectory inference. SpatialAnno is computationally efficient, easily scalable to spatially resolved transcriptomics with tens of thousands of spatial locations and thousands of genes (Supplementary Table S9). With simulation studies, we demonstrated that SpatialAnno presents improved spatial annotation accuracy with either correct, under- or over-specification of the number of cell/domain types, robustness to the marker gene misspecification and efficient leveraging of non-marker genes compared with other annotation methods.

We examined the SRT data generated using different platforms, such as 10x Visium, ST, Slide-seqV1/2 and seqFISH, with various spatial resolutions. Using both DLPFC 10x Visium datasets and mouse OB ST datasets with manual annotations, we demonstrated the improved annotation accuracy of SpatialAnno with the capability of recovering laminar structures, while the identified PAGA graph using embeddings in SpatialAnno recovers an almost linear trajectory from WM to layer 1. In DLPFC datasets, the domains identified were well matched with the elevated expression for marker genes, such as *PCP4* and *MOBP* that are marker genes for layer 5 and WM, respectively (23,37). Using mouse hippocampus Slide-seqV1/2 datasets, we demonstrated that SpatialAnno can successfully detect the primary hippocampal subfields for CA1, CA3 and DG, with almost a perfect correlation between cell-type proportions in both datasets and the elevated expression levels for *Wfs1*, *Cpne4* and *C1ql2* are well matched with CA1, CA3 and DG regions identified by SpatialAnno, respectively. *Wfs1* showed differential expression in hippocampal field CA1 and has been reported to be highly expressed in the CA1 region (49). *Cpne4*, a known marker gene for hippocampal subfield CA3, was highly expressed in a region identified as CA3 (50). In addition, *C1ql2*, a marker gene for dentate principal cells, was expressed in a region identified as DG (51). When applied to mouse embryo seqFISH datasets, SpatialAnno not only provided improved annotation accuracy, but uniquely estimated cell-type-aware embeddings leading to the identification of two trajectories in brain regions,



**Figure 5.** Spatial cell-type annotation of the mouse embryo dataset. **(A)** Bar plots of Kappa, mF1 and ACC showing the cell-type annotation accuracy of different methods. **(B)** Spatial annotations for ground truth, SpatialAnno, scSorter, SCINA, Garnett and CellAssign. **(C)** Top, expression levels of corresponding cell-type-specific marker genes. Bottom, annotations of ground truth and SpatialAnno are shown on each spot. **(D)** Left: latent time trajectory generated by slingshot on low dimensional embeddings of SpatialAnno. Right: clustering of the forebrain/midbrain/hindbrain cells into four spatially distinct clusters representing different regions of the developing brain.

originating in mesencephalon towards the rhombencephalon and prosencephalon, respectively. Moreover, cell-type distributions identified by SpatialAnno were well matched with the expression of their corresponding marker genes. For example, *Popdc2*, a cardiomyocyte marker, was expressed in the developing heart tube (52). *Foxa1*, a gut endoderm marker, showed the highest expression levels in the developing gut tube along the anterior-posterior axis of the embryo (53). In addition, *Foxf1*, a mesoderm marker that encodes a forkhead transcription factor expressed in the splanchnic mesenchyme surround-

ing the gut, was highly expressed at the identified splanchnic mesoderm (54).

SpatialAnno paves the way for future spatial annotation analyses in multiple scenarios. For example, a similar strategy can be applied to the problem of cell-type assignment in other spatial omics data, such as spatial resolved single-cell chromatin accessibility data (55) and spatial proteomics (56). To establish a complete spatial atlas of organism architecture, a critical bottleneck is to perform an automatic cell-type assignment with both considerations of molecular fea-

tures with/without prior knowledge as well as their spatial organization, SpatialAnno can substantially reduce both the irreproducibility and human effort in the processes of manual cell/domain-type assignment (56). We have primarily focused on examining Spatial Transcriptomics (SRT) technologies that measure high-dimensional gene expression at each tissue location. In addition, there are spatial proteomics technologies, such as Cytometry by Time-of-Flight (CyTOF) and CODEX, which characterize proteomic profiles of single cells using 30–40 protein channels (57). These technologies generate low-dimensional data. Since predefined marker proteins that define cell types are available, SpatialAnno can also be applied to analyze these datasets. In our analysis of real CyTOF data from breast cancer samples (58) and simulated CyTOF data from Cytomulate (unpublished manuscript), we observed that SpatialAnno and SCINA demonstrate similar performance. Furthermore, these methods outperform other approaches in terms of accuracy and robustness, as shown in Supplementary Figure S29.

The benefits of SpatialAnno come with some caveats that may require further exploration. First, SpatialAnno is applicable for spatial annotation in a single tissue slide. With multiple tissue slides available, methods that are capable of integrating multiple SRT datasets for cell/domain-type annotation are sincerely needed (59). Second, SpatialAnno was designed to perform annotation analysis of data with a single modality. However, incorporating multi-modal data with data of other modalities can further improve annotation accuracy. Third, many of the early SRT technologies do not have a single-cell resolution, and SpatialAnno is only able to assign domains with prior knowledge of each spot for those datasets. Cell-type annotation for this type of dataset further requires simultaneous deconvolution with spatial cellular annotation.

## Data availability

This study made use of publicly available datasets. These include the mouse OB dataset (<https://www.spatialresearch.org/>), DLPFC dataset on the 10x Visium platform are accessible at (<https://github.com/LieberInstitute/spatialLIBD>), seqFISH dataset (<https://doi.org/10.18129/B9.bioc.MouseGastrulationData>), and mouse hippocampus Slide-seq and Slide-seqV2 datasets ([https://singlecell.broadinstitute.org/single\\_cell/study/SCP948/robust-decomposition-of-cell-type-mixtures-in-spatial-transcriptomics](https://singlecell.broadinstitute.org/single_cell/study/SCP948/robust-decomposition-of-cell-type-mixtures-in-spatial-transcriptomics)). The SpatialAnno software and source code have been deposited at <https://github.com/Shufeyangyi2015310117/SpatialAnno>.

The code underlying this article is available in Zenodo at <https://doi.org/10.5281/zenodo.7414189>.

## Supplementary data

Supplementary Data are available at NAR Online.

## Acknowledgements

**Author contributions:** X.S. and J.L. initiated and designed the study. X.S. and Y.Y. developed the method, implemented the software, performed simulations and analyzed real data. X.S. and J.L. wrote the manuscript, and all authors edited and revised the manuscript.

## Funding

National Key R&D Program of China [2021YFA1000100, 2021YFA1000101]; University Development Fund from The Chinese University of Hong Kong, Shenzhen [UDF01003033]; National Natural Science Foundation of China [12171229, 71931004]; The Science and Technology Commission of Shanghai Municipality [22ZR1420500]. Funding for open access charge: The Science and Technology Commission of Shanghai Municipality [22ZR1420500].

## Conflict of interest statement

None declared.

## References

- Moffitt,J.R., Bambah-Mukku,D., Eichhorn,S.W., Vaughn,E., Shekhar,K., Perez,J.D., Rubinstein,N.D., Hao,J., Regev,A., Dulac,C., et al. (2018) Molecular, spatial, and functional single-cell profiling of the hypothalamic preoptic region. *Science*, **362**, eaau5324.
- Lohoff,T., Ghazanfar,S., Missarova,A., Koulena,N., Pierson,N., Griffiths,J., Bardot,E., Eng,C.-H., Tyser,R., Argelaguet,R., et al. (2022) Integration of spatial and single-cell transcriptomic data elucidates mouse organogenesis. *Nat. Biotechnol.*, **40**, 74–85.
- Stickels,R.R., Murray,E., Kumar,P., Li,J., Marshall,J.L., Di Bella,D.J., Arlotta,P., Macosko,E.Z. and Chen,F. (2021) Highly sensitive spatial transcriptomics at near-cellular resolution with Slide-seqV2. *Nat. Biotechnol.*, **39**, 313–319.
- Chen,A., Liao,S., Cheng,M., Ma,K., Wu,L., Lai,Y., Qiu,X., Yang,J., Xu,J., Hao,S., et al. (2022) Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA nanoball-patterned arrays. *Cell*, **185**, 1777–1792.
- Palla,G., Fischer,D.S., Regev,A. and Theis,F.J. (2022) Spatial components of molecular tissue biology. *Nat. Biotechnol.*, **40**, 308–318.
- Lähnemann,D., Köster,J., Szczurek,E., McCarthy,D.J., Hicks,S.C., Robinson,M.D., Vallejos,C.A., Campbell,K.R., Beerenwinkel,N., Mahfouz,A., et al. (2020) Eleven grand challenges in single-cell data science. *Genome Biol.*, **21**, 1–35.
- Franzén,O., Gan,L.-M. and Björkegren,J.L. (2019) PanglaoDB: a web server for exploration of mouse and human single-cell RNA sequencing data. *Database*, **2019**, Baz046.
- Zhang,X., Lan,Y., Xu,J., Quan,F., Zhao,E., Deng,C., Luo,T., Xu,L., Liao,G., Yan,M., et al. (2019) CellMarker: a manually curated resource of cell markers in human and mouse. *Nucleic Acids Res.*, **47**, D721–D728.
- Clarke,Z.A., Andrews,T.S., Atif,J., Pouyabahar,D., Innes,B.T., MacParland,S.A. and Bader,G.D. (2021) Tutorial: guidelines for annotating single-cell transcriptomic maps using automated and manual methods. *Nat. Protoc.*, **16**, 2749–2764.
- Duò,A., Robinson,M.D. and Soneson,C. (2018) A systematic performance evaluation of clustering methods for single-cell RNA-seq data. *F1000Research*, **7**, 1141.
- Kiselev,V.Y., Yiu,A. and Hemberg,M. (2018) scmap: projection of single-cell RNA-seq data across data sets. *Nat. Methods*, **15**, 359–362.
- Zhang,Z., Luo,D., Zhong,X., Choi,J.H., Ma,Y., Wang,S., Mahrt,E., Guo,W., Stawiski,E.W., Modrusan,Z., et al. (2019) SCINA: a semi-supervised subtyping algorithm of single cells and bulk samples. *Genes*, **10**, 531.
- Pliner,H.A., Shendure,J. and Trapnell,C. (2019) Supervised classification enables rapid annotation of cell atlases. *Nat. Methods*, **16**, 983–986.
- Zhang,A.W., O'Flanagan,C., Chavez,E.A., Lim,J.L., Ceglia,N., McPherson,A., Wiens,M., Walters,P., Chan,T., Hewitson,B., et al. (2019) Probabilistic cell-type assignment of single-cell RNA-seq

- for tumor microenvironment profiling. *Nat. Methods*, **16**, 1007–1015.
15. Guo,H. and Li,J. (2021) scSorter: assigning cells to known cell types according to marker genes. *Genome Biol.*, **22**, 1–18.
  16. Wu,F.-Y. (1982) The potts model. *Rev. Mod. Phys.*, **54**, 235.
  17. Zhang,Y., Brady,M. and Smith,S. (2001) Segmentation of brain MR images through a hidden Markov random field model and the expectation-maximization algorithm. *IEEE T. Med. Imaging*, **20**, 45–57.
  18. Zheng,G.X., Terry,J.M., Belgrader,P., Ryvkin,P., Bent,Z.W., Wilson,R., Ziraldo,S.B., Wheeler,T.D., McDermott,G.P., Zhu,J., et al. (2017) Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.*, **8**, 1–12.
  19. Butler,A., Hoffman,P., Smibert,P., Papalexi,E. and Satija,R. (2018) Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.*, **36**, 411–420.
  20. Hafemeister,C. and Satija,R. (2019) Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol.*, **20**, 1–15.
  21. Yang,Y., Shi,X., Liu,W., Zhou,Q., Chan Lau,M., Chun Tatt Lim,J., Sun,L., Ng,C.C.Y., Yeong,J. and Liu,J. (2022) SC-MEB: spatial clustering with hidden Markov random field using empirical Bayes. *Brief. Bioinform.*, **23**, bbab466.
  22. Liu,W., Liao,X., Yang,Y., Lin,H., Yeong,J., Zhou,X., Shi,X. and Liu,J. (2022) Joint dimension reduction and clustering analysis of single-cell RNA-seq and spatial transcriptomics data. *Nucleic Acids Res.*, **50**, e72.
  23. Maynard,K.R., Collado-Torres,L., Weber,L.M., Uytingco,C., Barry,B.K., Williams,S.R., Catallini,J.L., Tran,M.N., Besich,Z., Tippani,M. and et.al. (2021) Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex. *Nat. Neurosci.*, **24**, 425–436.
  24. Zhu,J., Sun,S. and Zhou,X. (2021) SPARK-X: non-parametric modeling enables scalable and robust detection of spatial expression patterns for large spatial transcriptomic studies. *Genome Biol.*, **22**, 1–25.
  25. Ståhl,P.L., Salmén,F., Vickovic,S., Lundmark,A., Navarro,J.F., Magnusson,J., Giacomello,S., Asp,M., Westholm,J.O., Huss,M., et al. (2016) Visualization and analysis of gene expression in tissue sections by spatial transcriptomics. *Science*, **353**, 78–82.
  26. Edsgård,D., Johnsson,P. and Sandberg,R. (2018) Identification of spatial expression trends in single-cell gene expression data. *Nat. Methods*, **15**, 339–342.
  27. Hao,Y., Hao,S., Andersen-Nissen,E., Mauck,W.M. III, Zheng,S., Butler,A., Lee,M.J., Wilk,A.J., Darby,C., Zager,M., et al. (2021) Integrated analysis of multimodal single-cell data. *Cell*, **184**, 3573–3587.
  28. Tepe,B., Hill,M.C., Pekarek,B.T., Hunt,P.J., Martin,T.J., Martin,J.F. and Arenkiel,B.R. (2018) Single-cell RNA-seq of mouse olfactory bulb reveals cellular heterogeneity and activity-dependent molecular census of adult-born neurons. *Cell Rep.*, **25**, 2689–2703.
  29. Saunders,A., Macosko,E.Z., Wysoker,A., Goldman,M., Krienen,F.M., de Rivera,H., Bien,E., Baum,M., Bortolin,L., Wang,S., et al. (2018) Molecular diversity and specializations among the cells of the adult mouse brain. *Cell*, **174**, 1015–1030.
  30. Berg,S., Kutra,D., Kroeger,T., Straehle,C.N., Kausler,B.X., Haubold,C., Schiegg,M., Ales,J., Beier,T., Rudy,M., et al. (2019) Ilastik: interactive machine learning for (bio) image analysis. *Nat. Methods*, **16**, 1226–1232.
  31. Lun,A.T., McCarthy,D.J. and Marioni,J.C. (2016) A step-by-step workflow for low-level analysis of single-cell RNA-seq data with Bioconductor. *F1000Research*, **5**, 2122.
  32. Chen,X., Chen,S., Song,S., Gao,Z., Hou,L., Zhang,X., Lv,H. and Jiang,R. (2022) Cell type annotation of single-cell chromatin accessibility data via supervised Bayesian embedding. *Nat. Mach. Intel.*, **4**, 116–126.
  33. Hubert,L. and Arabie,P. (1985) Comparing partitions. *J. classif.*, **2**, 193–218.
  34. Zappia,L., Phipson,B. and Oshlack,A. (2017) Splatter: simulation of single-cell RNA sequencing data. *Genome Biol.*, **18**, 1–15.
  35. Berger,J. and Pericchi,L. (2014) Bayes factors. *Wiley StatsRef: statistics reference online*. pp. 1–14.
  36. Kass,R.E. and Raftery,A.E. (1995) Bayes factors. *J. Am. Stat. Assoc.*, **90**, 773–795.
  37. Molyneaux,B.J., Arlotta,P., Menezes,J.R. and Macklis,J.D. (2007) Neuronal subtype specification in the cerebral cortex. *Nat. Rev. Neurosci.*, **8**, 427–437.
  38. Wolf,F.A., Hamey,F.K., Plass,M., Solana,J., Dahlin,J.S., Göttgens,B., Rajewsky,N., Simon,L. and Theis,F.J. (2019) PAGA: graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single cells. *Genome Biol.*, **20**, 1–9.
  39. Erwin,S.R., Sun,W., Copeland,M., Lindo,S., Spruston,N. and Cembrowski,M.S. (2020) A sparse, spatially biased subtype of mature granule cell dominates recruitment in hippocampal-associated behaviors. *Cell Rep.*, **31**, 107551.
  40. Nagai,Y., Sano,H. and Yokoi,M. (2005) Transgenic expression of Cre recombinase in mitral/tufted cells of the olfactory bulb. *genesis*, **43**, 12–16.
  41. van der Linden,C.J., Gupta,P., Bhuiya,A.I., Riddick,K.R., Hossain,K. and Santoro,S.W. (2020) Olfactory stimulation regulates the birth of neurons that express specific odorant receptors. *Cell Rep.*, **33**, 108210.
  42. Martín-López,E., Corona,R. and López-Mascaraque,L. (2012) Postnatal characterization of cells in the accessory olfactory bulb of wild type and reeler mice. *Front. Neuroanat.*, **6**, 15.
  43. Pijuan-Sala,B., Griffiths,J.A., Guibentif,C., Hiscock,T.W., Jawaid,W., Calero-Nieto,F.J., Mulas,C., Ibarra-Soria,X., Tyser,R.C., Ho,D.L.L., et al. (2019) A single-cell molecular map of mouse gastrulation and early organogenesis. *Nature*, **566**, 490–495.
  44. Simeone,A., Acampora,D., Gulisano,M., Stornaiuolo,A. and Boncinelli,E. (1992) Nested expression domains of four homeobox genes in developing rostral brain. *Nature*, **358**, 687–690.
  45. Bouillet,P., Chazaud,C., Oulad-Abdelghani,M., Dollé,P. and Chambon,P. (1995) Sequence and expression pattern of the Stra7 (Gbx-2) homeobox-containing gene induced by retinoic acid in P19 embryonal carcinoma cells. *Dev. Dynam.*, **204**, 372–382.
  46. Leimeister,C., Bach,A. and Gessler,M. (1998) Developmental expression patterns of mouse sFRP genes encoding members of the secreted frizzled related protein family. *Mech. Develop.*, **75**, 29–42.
  47. Wurst,W. and Bally-Cuif,L. (2001) Neural plate patterning: upstream and downstream of the isthmic organizer. *Nat. Rev. Neurosci.*, **2**, 99–108.
  48. Raible,F. and Brand,M. (2004) Divide et Impera—the midbrain–hindbrain boundary and its organizer. *Trends Neurosci.*, **27**, 727–734.
  49. Dong,H.-W., Swanson,L.W., Chen,L., Fanselow,M.S. and Toga,A.W. (2009) Genomic–anatomic evidence for distinct functional domains in hippocampal field CA1. *Proc. Natl. Acad. Sci.*, **106**, 11794–11799.
  50. Arneson,D., Zhang,G., Ying,Z., Zhuang,Y., Byun,H.R., Ahn,I.S., Gomez-Pinilla,F. and Yang,X. (2018) Single cell molecular alterations reveal target cells and pathways of concussive brain injury. *Nat. Commun.*, **9**, 1–18.
  51. Matsuda,K., Budisantoso,T., Mitakidis,N., Sugaya,Y., Miura,E., Kakegawa,W., Yamasaki,M., Konno,K., Uchigashima,M., Abe,M., et al. (2016) Transsynaptic modulation of kainate receptor functions by C1q-like proteins. *Neuron*, **90**, 752–767.
  52. Breher,S.S., Mavridou,E., Brenneis,C., Froese,A., Arnold,H.-H. and Brand,T. (2004) Popeye domain containing gene 2 (Popdc2) is a myocyte-specific differentiation marker during chick heart development. *Development. Dynam.*, **229**, 695–702.
  53. Sasaki,H. and Hogan,B. (1993) Differential expression of multiple fork head related genes during gastrulation and axial pattern formation in the mouse embryo. *Development*, **118**, 47–59.
  54. Mahlapuu,M., Ormestad,M., Enerback,S. and Carlsson,P. (2001) The forkhead transcription factor Foxf1 is required for

- differentiation of extra-embryonic and lateral plate mesoderm. *Development*, **128**, 155–166.
55. Deng,Y., Bartosovic,M., Ma,S., Zhang,D., Kukanja,P., Xiao,Y., Su,G., Liu,Y., Qin,X., Rosoklija,G.B., *et al.* (2022) Spatial profiling of chromatin accessibility in mouse and human tissues. *Nature*, **609**, 375–383.
56. Brbić,M., Cao,K., Hickey,J.W., Tan,Y., Snyder,M.P., Nolan,G.P. and Leskovec,J. (2022) Annotation of spatially resolved single-cell data with STELLAR. *Nat. Methods*, **19**, 1411–1418.
57. Wang,K., Yang,Y., Wu,F., Song,B., Wang,X. and Wang,T. (2023) Comparative analysis of dimension reduction methods for cytometry by time-of-flight data. *Nat. Commun.*, **14**, 1836.
58. Ali,H.R., Jackson,H.W., Zanotelli,V.R., Danenberg,E., Fischer,J.R., Bardwell,H., Provenzano,E., Rueda,O.M., Chin,S.-F., *et al.* (2020) Imaging mass cytometry and multiplatform genomics define the phenogenomic landscape of breast cancer. *Nat. Cancer*, **1**, 163–175.
59. Liu,W., Liao,X., Luo,Z., Yang,Y., Lau,M.C., Jiao,Y., Shi,X., Zhai,W., Ji,H., Yeong,J., *et al.* (2023) Probabilistic embedding, clustering, and alignment for integrating spatial transcriptomics data with PRECAST. *Nat. Commun.*, **14**, 296.